

**XAI for
Dummies**



Our Study

Ulrike Kuhl

Bielefeld University, Bielefeld, Germany

Mark T. Keane

University College Dublin, Dublin, Ireland

Six Steps to Heaven...



Motivation



Design



Materials & Procedure



Piloting



Data Collection & Analysis



Results





To what end do you want to test
what with whom in which context?



To what
end?

Examine how the **personal relevance** of AI system decisions impacts user behavior and perception of explanations.

What?

Human **prediction accuracy** w.r.t. system decisions that are:

- **either relevant for oneself or another person**, and
 - either supported by **counterfactuals or control** descriptions
- + subjective satisfaction & trust in the corresponding system

Whom?

AI experts attending this tutorial @IJCAI 2024



To what end do you want to test
what with whom in which context?



Context?

AI-based performance evaluation and
reward system *PERS*

Factor 1: self-relevant



Case 1/16

Please review the information and take a guess.

YOUR performance:

Productivity	Low
Quality of work	Low
Customer feedback	High
Punctuality	High
Bonus?	??

Will YOU get the bonus?

vs.

other-relevant



Case 1/16

Please review the information and take a guess.

Frank's performance:

Punctuality	High
Quality of work	Low
Productivity	High
Customer feedback	Low
Bonus?	??

Will Frank get the bonus?



To what end do you want to test
what with whom in which context?



Context?

AI-based performance evaluation and
reward system *PERS*

Factor 2: counterfactual vs. control description

Bonus? ☒ Yes

You predicted: ☒ No
The system predicted: ☒ Yes

Explanation: If "Quality of work" of Emily's work had been "Low", then Emily would have NOT received the bonus.

[Next case](#)

Bonus? ☒ Yes

You predicted: ☒ No
The system predicted: ☒ Yes

Explanation: If "Quality of work" of YOUR work had been "Low", then YOU would have NOT received the bonus.

[Next case](#)

Bonus? ☒ Yes

You predicted: ☒ No
The system predicted: ☒ Yes

Explanation: The system predicted Yes because of the aspects of Laura's performance.

[Next case](#)

Bonus? ☒ Yes

You predicted: ☒ No
The system predicted: ☒ Yes

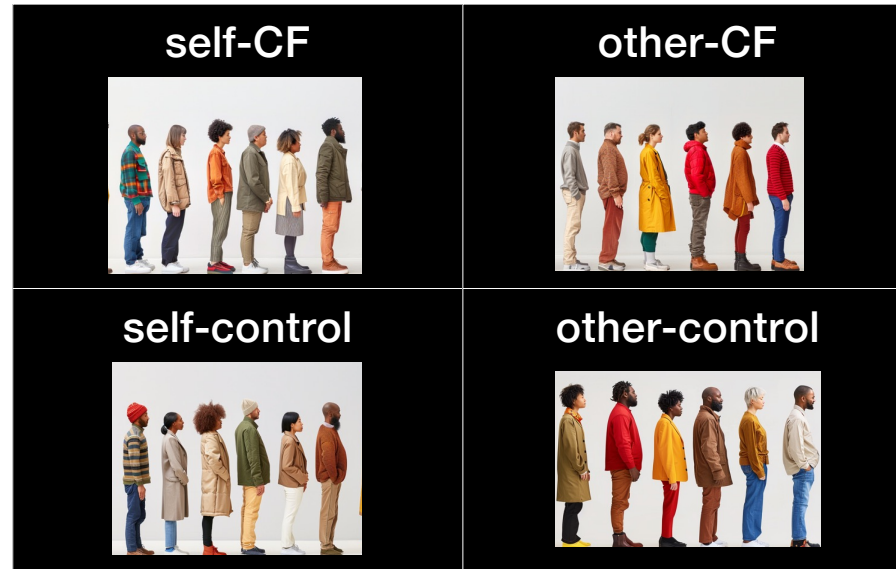
Explanation: The system predicted Yes because of the aspects of YOUR performance.

[Next case](#)



Design-E

2 x 2



We used a ***between-participant*** design, with ***two factors*** (self-relevance, explanation), with **two levels** each:
self-relevance has levels **self** vs. **other**
explanation has levels **counterfactual** vs. **control**



material-set

Back to balancing and randomizing



Materials &
Procedure

Quality of work	H	H	H	H	H	H	H	H	L	L	L	L	L	L	L	L
Productivity	H	H	H	H	L	L	L	L	H	H	H	H	L	L	L	L
Punctuality	H	H	L	L	H	H	L	L	H	H	L	L	H	H	L	L
Customer feedback	H	L	H	L	H	L	H	L	H	L	H	L	H	L	H	L
Bonus?	Y	Y	Y	Y	Y	Y	N	N	Y	Y	N	N	N	N	N	N

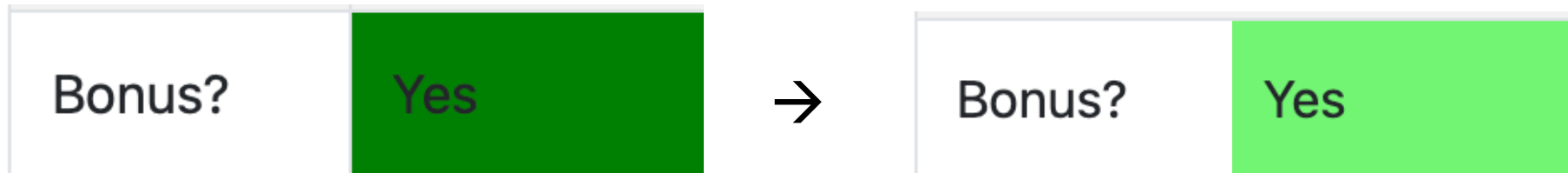
- ✓ 1/2 items are positive 1/2 items are negative
- ✓ presented order of features randomized across participants
- ✓ presented order of items randomized across participants
- ✓ other-relevant condition: 1/2 items shown as male gendered, 1/2 items shown as female gendered



Pilot: Done with family and friends beforehand



N=37 → 1 failed BOTH attention checks → N=36



Instruction "other" conditions:

Imagine SEVERAL EMPLOYEES THAT work...

→ Imagine SEVERAL EMPLOYEES that work...

Summary of the number of participants in each group:

	cond	Number of Participants
0	o-cfe	16
1	o-con	8
2	s-cfe	7
3	s-con	5

Randomization produced
unbalanced groups!



Piloting

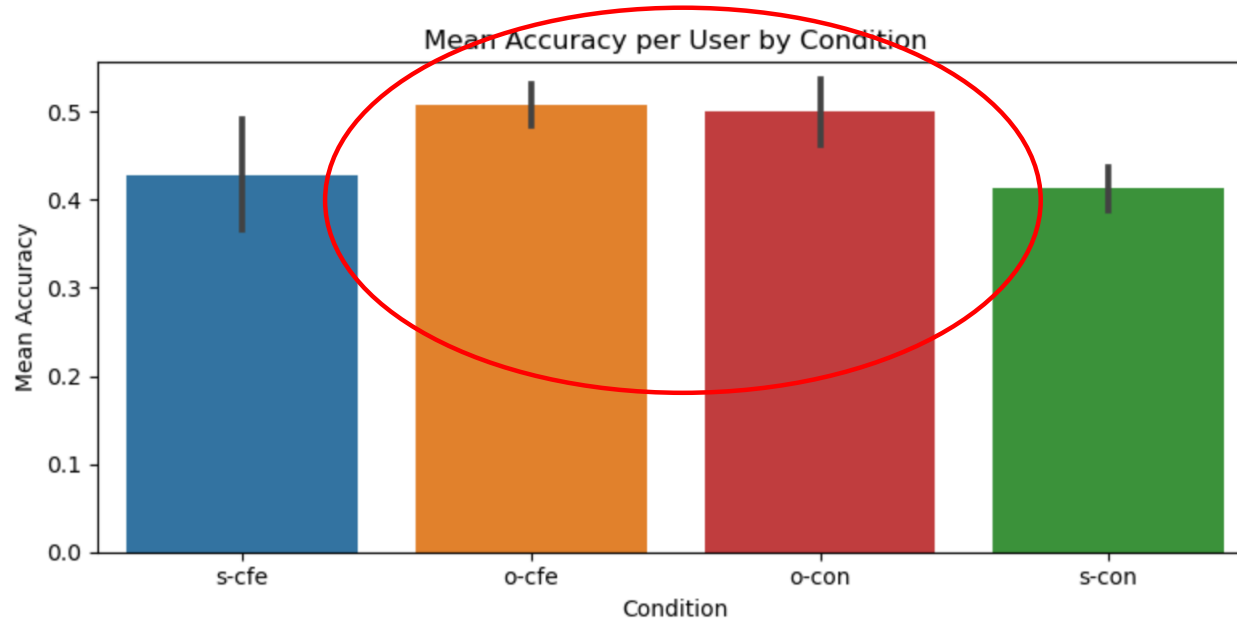
Pilot: Done with family and friends beforehand

Mean accuracy over task:

Punctuality	High
Bonus?	??

Will YOU get the bonus?

☐ Yes ☐ No



looks like 'other' > 'self'! (ANOVA confirmed!)



Pilot: Done with family and friends beforehand



Piloting

How close to ground truth?:
looks like CFE > CON

Please rate your agreement with the following statement as truthfully as possible.

From your interaction with the system, please weight the four aspects used by PERS according to their importance

Quality of work

Least Important

Most important

Productivity

Least Important

Most important

Punctuality

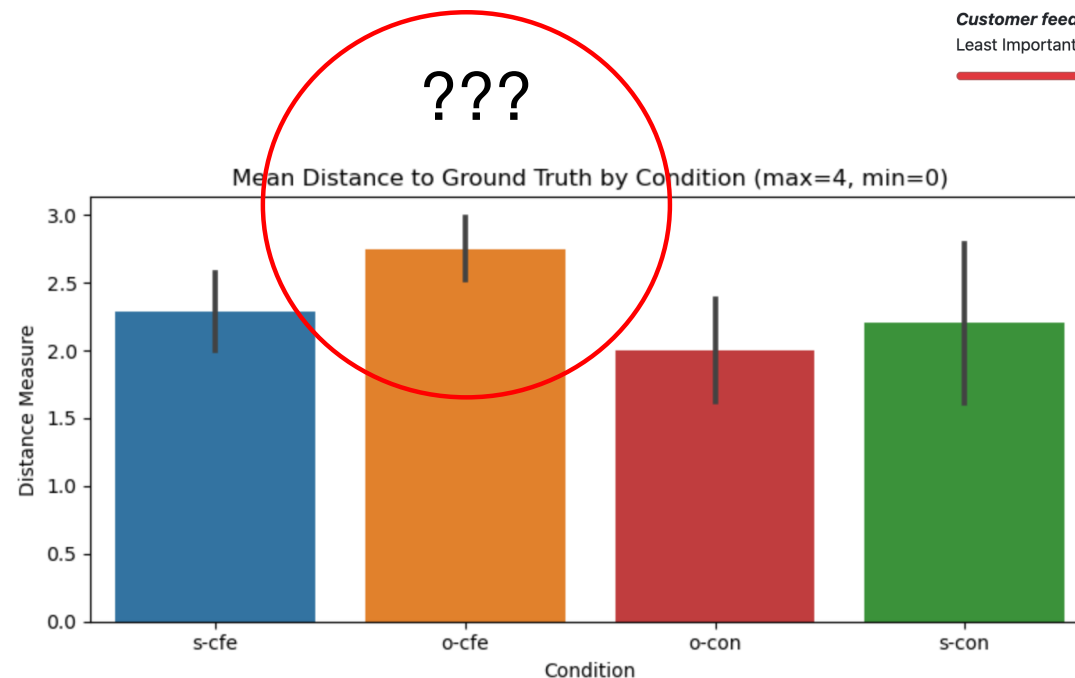
Least Important

Most important

Customer feedback

Least Important

Most important



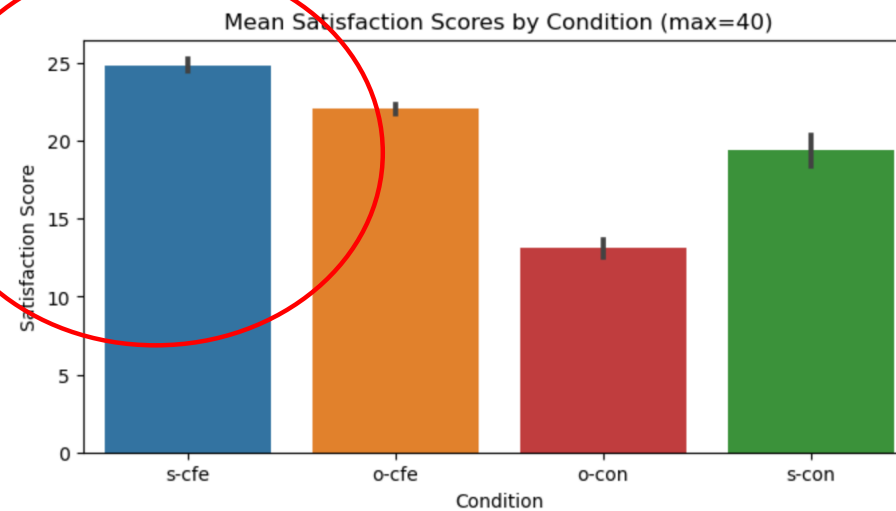


Pilot: Done with family and friends beforehand

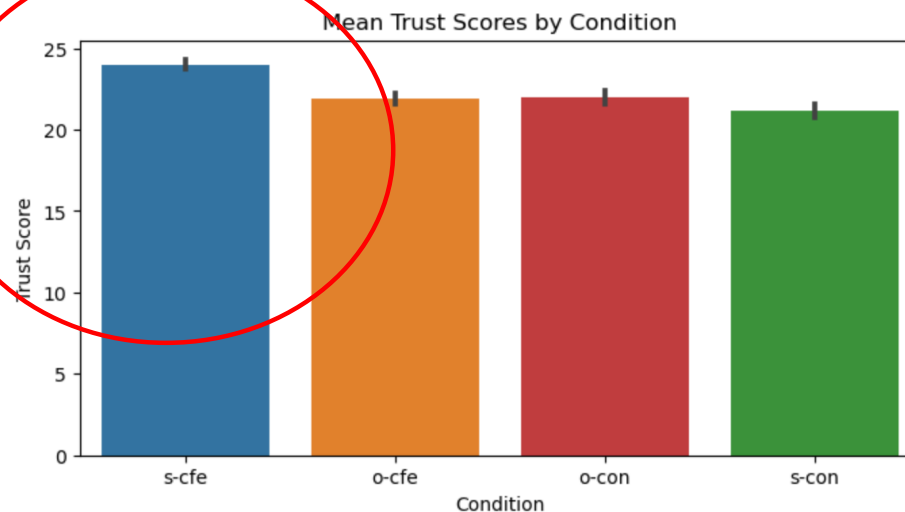


Piloting

Satisfaction



Trust



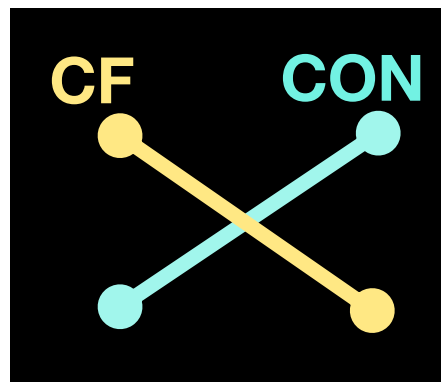
“s-cfe”
higher
than all
others?

(Before) Starting the Statistics



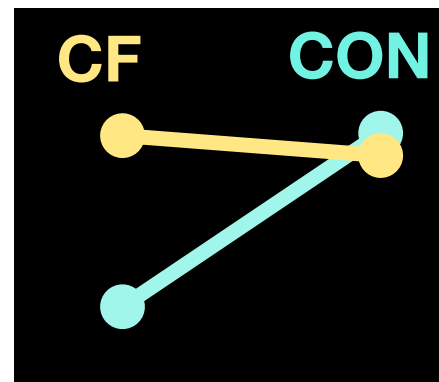
Data Collection
& Analysis

→ A quick look at 2x2 designs, and how the means could behave



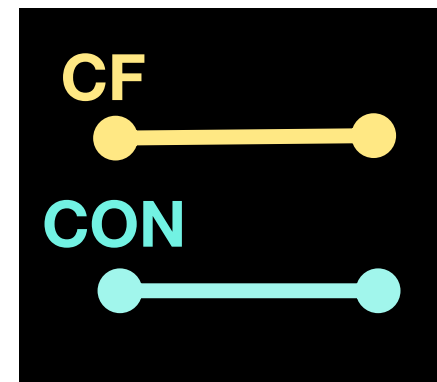
self other

*strong
interaction*



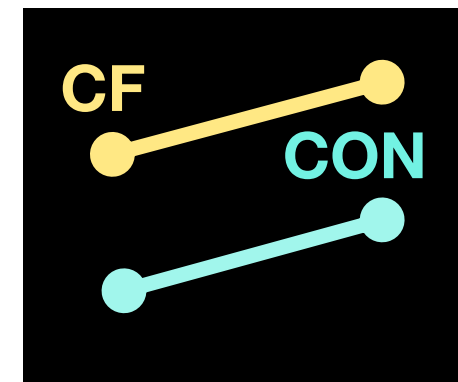
self other

*weaker
interaction*



self other

*Independent,
no difference*



self other

*independent,
with difference*

Doing the Statistics – for real!



**Data Collection
& Analysis**



Results



Let's do it live!

https://colab.research.google.com/github/andreArtelt/IJCAI24-CF_Tut/blob/main/CESORP_study_evaluation.ipynb

